

Solution of Banff Challenge 2 - Wolfgang Rolke

My solution for both problems is based on the likelihood ratio test statistic

$$\lambda(\mathbf{x}) = 2 \left(\max \{ \log L(\theta | \mathbf{x}) : \theta \} - \max \{ \log L(\theta | \mathbf{x}) : \theta \in \Theta^0 \} \right)$$

According to standard theorems in Statistics $\lambda(\mathbf{X})$ often has a χ^2 distribution with the number of degrees of freedom the difference between the number of free parameters and the number of free parameters under the null hypothesis. This turns out to be true for problem 2 but not for problem 1, in which case the null distribution can be found via simulation.

Problem 1

Here we have:

$$f(x) = 0.9999546e^{-10x}, \quad 0 < x < 1$$

$$\varphi(x; E) = \frac{1}{\sqrt{2\pi} \cdot 0.03} e^{-\frac{1}{2} \frac{(x-E)^2}{0.03^2}}$$

$$g(x; E) = \frac{\varphi(x; E)}{\int_0^1 \varphi(t; E) dt}, \quad 0 < x < 1$$

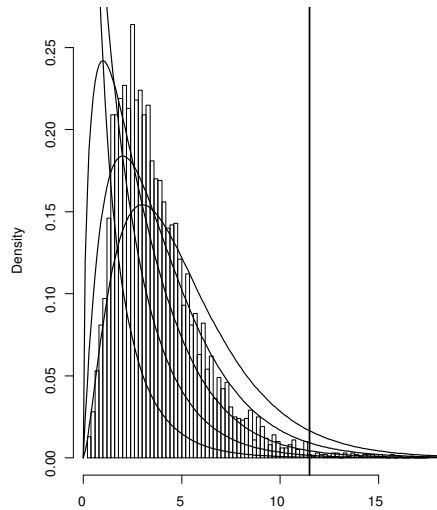
$$h(x; \alpha, E) = (1 - \alpha)f(x) + \alpha g(x; E)$$

$$H_0 : \alpha = 0 \text{ vs } H_a : \alpha > 0$$

$$\log L(\alpha, E | \mathbf{x}) = \sum_{i=1}^n \log[(1 - \alpha)f(x_i) + \alpha g(x_i; E)]$$

Now $\max \{ \log L(\alpha, E | \mathbf{x}) \}$ is the loglikelihood function evaluated at the maximum likelihood estimator and $\max \{ \log L(\alpha, E | \mathbf{x}) : \theta \in \Theta^0 \} = \log L(0, 0 | \mathbf{x})$. Note that if $\alpha = 0$ any choice of E yields the same value of the likelihood function.

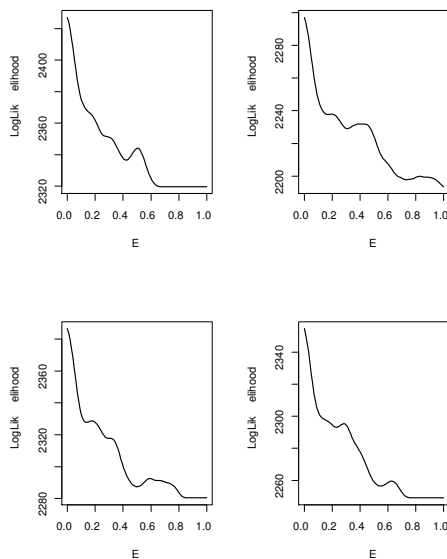
In the following figure we have the histogram of 100000 values of $\lambda(\mathbf{x})$ for a simulation with $n = 500$ and $\alpha = 0$ together with the densities of the χ^2 distribution with df's from 1 to 5. Clearly none of these yields an acceptable fit. Instead we use the simulated data to find the 99% quantile and reject the null hypothesis if $\lambda(\mathbf{x})$ is larger than that, shown as the vertical line in the graph.



In general the critical value will depend on the sample size, but for those in the challenge (500-1500) it is always about 11.5.

If it was decided to do discovery using 5σ the critical value can be found using importance sampling.

Finding the mle is a non-trivial exercise because there are many local minima. The next figure shows the log-likelihood as a function of E with α fixed at 0.05 for 4 cases.



To find the mle i used a two-step procedure: first a fine grid search over values of E from -0.015 to 1 in steps of 0.005. At each value of E the corresponding value of α that maximizes the log-likelihood is found. In a second step i start at the best point found above

and use Newton-Raphson to find the overall mle.

Problem 2

Again i want to use:

$$h(x; \alpha, \beta) = (1 - \alpha - \beta)f_1(x) + \beta f_2(x) + \alpha g(x)$$

$$H_0 : \alpha = 0 \text{ vs } H_a : \alpha > 0$$

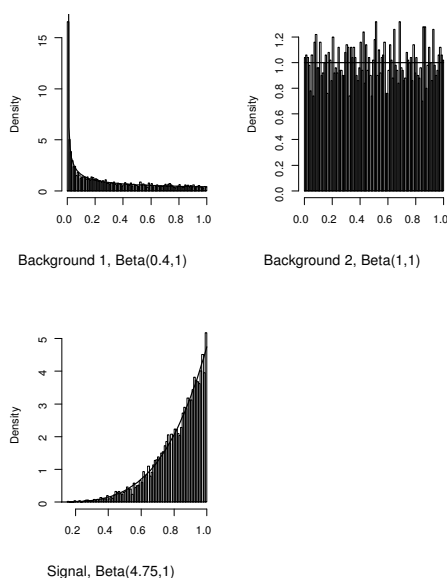
$$\log L(\alpha, \beta | \mathbf{x}) = \sum_{i=1}^n \log[(1 - \alpha - \beta)f_1(x) + \beta f_2(x) + \alpha g(x)]$$

Now $\max\{\log L(\alpha, \beta | \mathbf{x})\}$ is the loglikelihood function evaluated at the maximum likelihood estimator and $\max\{\log L(\alpha, \beta | \mathbf{x}) : \theta \in \Theta^0\} = \max\{\log L(0, \beta | \mathbf{x}) : \beta\}$.

The difficulty is of course that we don't know f_1 , f_2 or g . I have used three different ways to find them:

Parametric Fitting

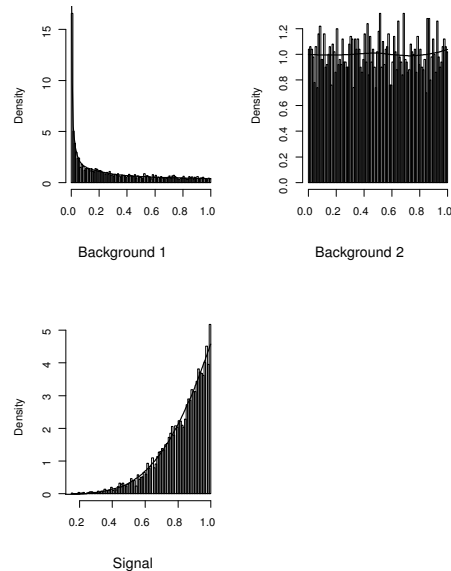
Here one tries to find a parametric density that gives a reasonable fit to the data. For the data in the challenge this turns out to be very easy. In all three cases a Beta density gives a very good fit:



Nonparametric Fitting:

There are a variety of methods known in Statistics for non-parametric density estimation. The difficulty with the data in the challenge is that it is bounded on a finite interval, a very common feature in HEP data. Moreover the slope of the density of Background 1 at 0 is infinite. I checked a number of methods and eventually ended up using the following: for Background 2, the Signal and the right half of Background 1 i bin the data (250 bins) find the counts and scale them to integrate to unity. Then i use the non-parametric density estimator loess from R with the default span (smoothing parameter). This works well except on the left side of Background 1. There the infinite slope of the

density would require a smoothing parameter that goes to 0. Instead i transform the data with $\log \frac{x}{1-x}$. The resulting data has a density without boundary, which i estimate using the routine density from R, again with the default bandwidth. This is then back-transformed to the 0-1 scale. This works well for the left side but not the right one, and so i "splice" the two densities together in the middle. The resulting densities are shown here:



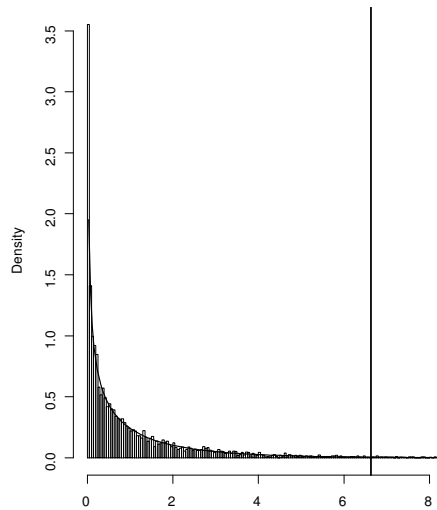
Semiparametric Fitting

It is possible to combine the two approaches above: fit some of the data parametrically and others non-parametrically, for example if the signal is known to have a Gaussian distribution but the background density is Monte Carlo data.

For the data in the challenge the three methods give very similar results, so i am submitting only the solution using the parametric fits.

Back to the Test

What is the null distribution of $\lambda(\mathbf{x})$ now? In the following figure we have the histogram of 5000 values of $\lambda(\mathbf{x})$ for a simulation with 500 events from Background 1, 100 events from Background 2 and no Signal events. $\alpha = 0$ together with the density of a χ^2 distribution with 1 df. The densities are fit parametrically.



Clearly this yield a very good fit, so we will reject the null hypothesis if $\lambda(\mathbf{x}) > q\chi^2(0.99, 1) = 6.635$, the 99th percentile of a chi-square distribution with 1 degree of freedom. The same result holds if the fitting was done non parametrically or semi parametrically.

Error Estimation

We have the following large sample theorem for maximum likelihood estimators: if we have a sample X_1, \dots, X_n with density $f(x; p)$ and \hat{p} is the mle for the parameter p , then under some regularity conditions

$$\sqrt{n} (\hat{p} - p) \sim N(0, \sigma)$$

where $\sigma^2 = -1/nE[\frac{d^2}{dp^2} \log f(X; p)]$, the Fisher information number. This can be estimated using the observed Fisher information number $\frac{1}{n} \sum_{i=1}^n \frac{d^2}{dp^2} \log h(X_i; p)$.

For problem 1 we find:

$$\text{Signal density: } S(x; E) = \frac{g(x; E)}{q(E)}$$

$$g(x; E) = dn(x; E, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2} \frac{(x-E)^2}{\sigma^2}}$$

$$pn(x; E, \sigma) = \int_{-\infty}^x dn(t; E, \sigma) dt$$

$$\varphi(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2}$$

$$q(E) = 1 - pn(0; E, \sigma) = 1 - pn(-E/\sigma; 0, 0)$$

$$g_E = \frac{dg}{dE} = g \frac{x-E}{\sigma^2}$$

$$g_{EE} = \frac{d^2g}{dE^2} = g \left(\frac{x-E}{\sigma^2} \right)^2 + g \frac{-1}{\sigma^2} = g \left[\left(\frac{x-E}{\sigma^2} \right)^2 - \frac{1}{\sigma^2} \right]$$

$$\frac{d}{dE} q(E) = -\varphi(-E/\sigma) \left(-\frac{1}{\sigma} \right) = \varphi(E/\sigma)/\sigma$$

$$\frac{d}{dx} \varphi(x) = -x\varphi(x)$$

$$r(E) = \log S(x; E) = \log g - \log q = \text{const} - \frac{(x-E)^2}{2\sigma^2} - \log q$$

$$r' = \frac{x-E}{\sigma^2} - \frac{\varphi(E/\sigma)/\sigma}{q}$$

$$r'' = -\frac{1}{\sigma^2} + \frac{-\varphi(E/\sigma) \frac{E}{\sigma^2} q - \varphi(E/\sigma)^2/\sigma^2}{q^2} =$$

$$-\frac{1}{\sigma^2} \left[1 + \varphi(E/\sigma) \frac{qE + \varphi(E/\sigma)}{q^2} \right]$$

$$\frac{dS}{dE} = S r' = S \left[\frac{x-E}{\sigma^2} - \frac{\varphi(E/\sigma)/\sigma}{q} \right]$$

$$\frac{d^2S}{dE^2} = S(r'' + r'^2) =$$

$$S \left[-\frac{1}{\sigma^2} \left[1 + \varphi(E/\sigma) \frac{E[1-\Phi(0, E)] + \varphi(E/\sigma)}{[1-\Phi(0, E)]^2} \right] + \left(\frac{x-E}{\sigma^2} + \frac{\varphi(E/\sigma)/\sigma}{1-\Phi(0, E)} \right)^2 \right]$$

$$\psi(x; \alpha, E) = (1 - \alpha)f(x) + \alpha S(x; E)$$

$$\sum_{i=1}^n \log \psi(x_i; \alpha, E) =$$

$$\sum_{i=1}^n \log[(1 - \alpha)f(x_i) + \alpha S(x_i; E)]$$

$$\frac{d\psi}{d\alpha} = S - f$$

$$\frac{d\psi}{dE} = \alpha S' = \alpha S \left[\frac{x-E}{\sigma^2} - \frac{\varphi(E/\sigma)/\sigma}{q} \right]$$

$$\frac{d \log \psi}{d\alpha} = \frac{\psi_\alpha}{\psi} = \frac{S-f}{\psi}$$

$$\frac{d \log \psi}{dE} = \frac{\psi_E}{\psi} = \alpha \frac{S}{\psi} \left[\frac{x-E}{\sigma^2} - \frac{\varphi(E/\sigma)/\sigma}{q} \right]$$

$$\frac{d^2 \log \psi}{d\alpha^2} = -\frac{(S-f)^2}{\psi^2}$$

$$\frac{d^2 \log \psi}{d\alpha dE} = \frac{S' \psi - (S-f) \psi_E}{\psi^2}$$

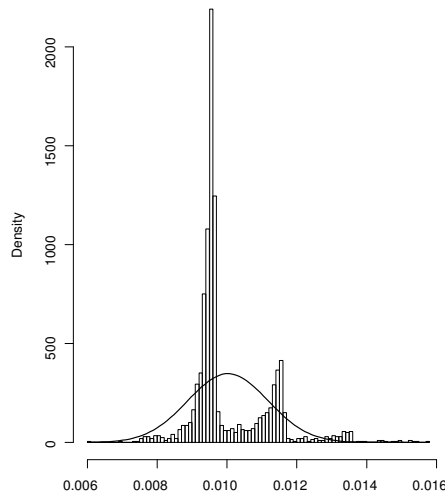
$$\frac{d^2 \log \psi}{dE^2} = \frac{d}{dE} \left[\frac{\alpha S'}{\psi} \right] = \frac{\alpha S'' \psi - \psi_g^2}{\psi^2}$$

so the standard error of α is estimated with $\left[\sum_{i=1}^n \left(\frac{S(X_i; \hat{\alpha}, \hat{E}) - f(X_i)}{\psi(X_i; \hat{\alpha}, \hat{E})} \right)^2 \right]^{-1/2}$ and the standard

error of E is given by $\left[\sum_{i=1}^n \left(\frac{\alpha S'' \psi - \psi_g^2}{\psi^2} \right)^2 \right]^{-1/2}$.

For problem 2 the errors are found similarly.

How good are these errors? As always with large sample theory, there is a question whether it works for a specific problem at the available sample sizes. Here is the result of a mini MC study: we generate 500 events from the background of problem 1 and 4 events from the signal with $E = 0.8$. This is repeated 2000 times. The histogram of estimates for the mixing ratio α is shown here:



Clearly the distribution of the estimates is not Gaussian, and therefore the errors are wrong.

So we need a different method for finding the 68% intervals. This can be done via the statistical bootstrap. In our submission we will include intervals based on the bootstrap, specifically the 16th and the 84th percentile of a bootstrap sample of size 300. In a real live problem it would be easy to run a mini MC for a specific case and then decide which errors to use. Because we have to process 20000 data sets we can not do so for each individually, and therefore use of the bootstrap intervals.

Power Studies

Problem 1:

I used the following code to do the power study:

```
counter=0;
for(int k=0;k<10000;++k) {
    nsig=rpois(0.07519885*D,seed);
    for (int i=0; i<nsig; ++i) x[i]=rnorm(E,seed);
    nback=rpois(1000,seed);
    for (int i=0; i<nback; ++i) x[nsig+i]=rbackground(seed)
    lrt=findLRT(x);
```

```

    if(lrt>11.5) ++counter;
}
power=counter/M;

```

The results are :

(D, E)	Power
(1010, 0.1)	0.356
(137, 0; 0.5)	0.457
(18, 0; 0, 9)	0.184

Problem 2:

I used the following code to do the power study:

```

counter=0;
for(int k=0;k<10000;++k) {
    n1=rnorm(1,900,90)
    x1=sample(bc2p2b1mc,size=n1)
    n2=max(rnorm(1,100,100),0)
    x2=sample(bc2p2b2mc,size=n2)
    n3=rpois(75,seed);
    x3=sample(bc2p2sigmc,size=n3)
    nback=rpois(1000,seed);
    x=c(x1,x2,x3)
    lrt=findLRT(x);
    if(lrt>6.635) ++counter;
}
power=counter/M;

```

The result is a power of **88%**.

Submission

For problem 1 i attach the file Rolkep1.dat. The first lines are:

```

Num Signal p-val SigLoc 68%Low 68%High D68%Low D68%High
0 0 0.9017 0 0 0 0 0
1 0 0.6676 0 0 0 0 0
2 1 0.0011 0.381 0.365 0.403 207.9 382.4
3 0 0.8895 0 0 0 0 0

```

so the first dataset for which i claim a signal is #2 (the third), with a p-value of 0.001, a signal at 0.381, 68% CI (0.365, 0.403) and a 68% CI for D of (207.9, 382.4)

Note that in terms of the number of signal events this means (15.6, 28.8).

For problem 2 i attach the file Rolkep2.dat. The first lines are:

Num Signal p-val NumSigEvents 68%Low 68%High

0 0 0.3937 0 0 0

1 1 5e-04 76 57.6 97.9

2 0 0.5942 0 0 0

so i claim a signal for dataset #1 with a p-value of 0.0005, 76 signal events with a 68% CI of (57.6, 97.9)